

COMPSCI 690RA: Randomized Algorithms and Probabilistic Data Analysis

Prof. Cameron Musco

University of Massachusetts Amherst. Spring 2022.

Lecture 11

Logistics

- Problem Set 4 is due next Tuesday 5/3 at 8pm.
- The final exam is next Friday 5/5 at 10:30am for those that are taking it.
- I will hold extended office hours Wed. 5/3 from 2-4pm and Thurs. 5/4 from 4-6pm.
- I will accept final project submissions up until Sunday 5/8 at 11:59pm.
- SRTI's are open for this course. It would be very helpful to me if you can fill them out!
- This was my first time teaching this course, so feedback on what worked and what didn't is really useful to me.

Last Week: More Advanced Markov Chains.

- The gambler's ruin problem.
- Start on Markov chains and their analysis.
- Aperiodicity and stationary distribution of a Markov chain.
- Start on mixing time.

Today: Finish up Markov Chains

- Mixing time analysis via coupling.
- Markov Chain Monte Carlo (MCMC) methods.

Fundamental Theorem of Markov Chains

Theorem (The Fundamental Theorem of Markov Chains)

Let X_0, X_1, \dots be a Markov chain with a finite state space and transition matrix $P \in [0, 1]^{m \times m}$. If the chain is both irreducible and aperiodic,

1. There exists a **unique** stationary distribution $\pi \in [0, 1]^m$ with $\pi = \pi P$.
2. For any states i, j , $\lim_{t \rightarrow \infty} \Pr[X_t = i | X_0 = j] = \pi(i)$. I.e., for any initial distribution q_0 , $\lim_{t \rightarrow \infty} q_t = \lim_{t \rightarrow \infty} q_0 P^t = \pi$.

Question for today: How long does it take us to converge close to this stationary distribution?

Mixing Time

Definition (Mixing Time)

Consider a Markov chain $\mathbf{X}_0, \mathbf{X}_1, \dots$ with unique stationary distribution π . Let $q_{i,t}$ be the distribution over states at time t assuming $\mathbf{X}_0 = i$. The mixing time is defined as:

$$\tau(\epsilon) = \min \left\{ t : \max_{i \in [m]} \|q_{i,t} - \pi\|_{TV} \leq \epsilon \right\}.$$

Note: If $\|q_{i,t} - \pi\|_{TV} \leq \epsilon$ then for any $t' \geq t$, $\|q_{i,t'} - \pi\|_{TV} \leq \epsilon$. **Coupling Motivation:** Last time we showed that

$$\max_{i \in [m]} \|q_{i,t} - \pi\|_{TV} \leq \max_{i,j \in [m]} \|q_{i,t} - q_{j,t}\|_{TV}.$$

By **Kantorovich-Rubinstein duality**, for $\mathbf{X}_t, \mathbf{Y}_t$ distributed by evolving the chain for t steps starting from state i or j respectively, we have:

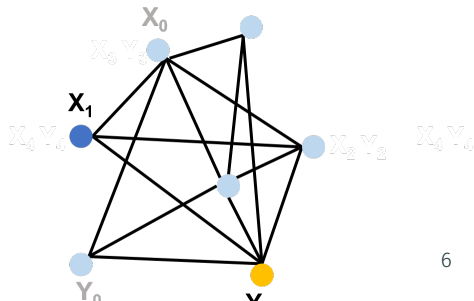
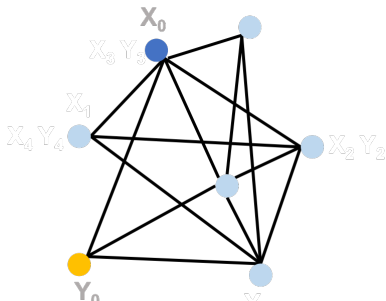
$$\max_{i,j \in [m]} \|q_{i,t} - q_{j,t}\|_{TV} \leq \max_{i,j \in [m]} \Pr[\mathbf{X}_t \neq \mathbf{Y}_t].$$

Formal Coupling Definition

Definition (Coupling)

For a finite Markov chain X_0, X_1, \dots with transition matrix $P \in \mathbb{R}^{m \times m}$, a coupling is a joint process $(X_0, Y_0), (X_1, Y_1), \dots$ such that:

1. $X_0 = i$ and $Y_0 = j$ for some $i, j \in [m]$.
2. $\Pr[X_t = j | X_{t-1} = i] = \Pr[Y_t = j | Y_{t-1} = i] = P_{i,j}$
3. If $X_t = Y_t$, then $X_{t+1} = Y_{t+1}$.



Coupling Example: Mixing Time of Shuffling

How many times do we need to swap a random card to the top of the deck so that the distribution of orderings on our cards is ϵ -close in TV distance to the uniform distribution over all permutations?

Coupling:

- Let X_0, X_1, \dots be the Markov chain where a random card is moved to the top in each step.
- Let Y_0, Y_1 be a correlated Markov chain. When card S is swapped to the top in the X chain, swap S to the top in the Y chain as well.
- Can check that this is a valid coupling since X_t, Y_t have the correct marginal distributions, and since
$$X_t = Y_t \implies X_{t+1} = Y_{t+1}$$
- Observe that $X_t = Y_t$ as soon as all c unique cards have been swapped at least once. How many swaps does this take?

Coupling Example: Mixing Time of Shuffling

$$\begin{aligned}\max_{i \in [m]} \|q_{i,t} - \pi\|_{TV} &\leq \max_{i,j \in [m]} \Pr[\mathbf{T}_{i,j} > t] \\ &\leq \Pr[\text{< } c \text{ unique cards are swapped in } t \text{ swaps}]\end{aligned}$$

By coupon collector analysis for $t \geq c \ln(c/\epsilon)$, this probability is bounded by ϵ . In particular, by the fact that $(1 - \frac{1}{c})^{c \ln c/\epsilon} \leq \frac{\epsilon}{c}$ plus a union bound over c cards.

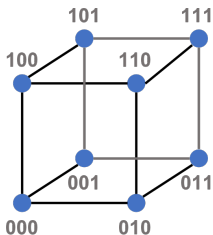
Thus, for $t \geq c \ln(c/\epsilon)$,

$$\max_{i \in [m]} \|q_{i,t} - \pi\|_{TV} \leq \max_{i,j \in [m]} \|q_{i,t} - q_{j,t}\|_{TV} \leq \epsilon.$$

I.e., $\tau(\epsilon) \leq c \ln(c/\epsilon)$.

Coupling Example: Random Walk on a Hypercube

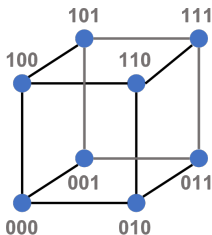
Let X_0, X_1 be a Markov chain over state space $\{0, 1\}^n$. In each step, pick a random position $i \in [n]$ and set $X_t(i) = 0$ with probability $1/2$ and $X_t(i) = 1$ with probability $1/2$.



What is a coupling $(X_0, Y_0), (X_1, Y_1), \dots$ on this chain that we can use to bound the mixing time of this walk?

Coupling Example: Random Walk on a Hypercube

In each step, pick a single random position $i \in [n]$ and let $X_t(i) = Y_t(i) = 0$ with probability $1/2$ and $X_t(i) = Y_t(i) = 1$ with probability $1/2$.



How large must we set t so that $\Pr[X_t \neq Y_t] \leq \epsilon$?

Upshot: The mixing time of the n -dimensional hypercube is $\tau(\epsilon) = O(n \log(n/\epsilon))$.

Markov Chain Monte Carlo

Markov Chain Monte Carlo

Many applications in computational biology, machine learning, theoretical computer science, etc. require sampling from complex distributions, which are difficult to write down in closed form, and difficult to directly sample from.

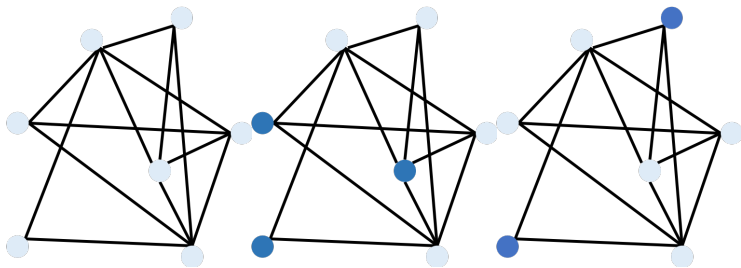
A very common approach is to design a Markov chain whose **stationary distribution π is equal to the distribution of interest.**

By running this Markov chain for at least $\tau(\epsilon)$ steps (burn-in time), one can draw a sample which is nearly from the distribution of interest.

Note: A major focus is on designing and analyzing Markov chains where $\tau(\epsilon)$ is small. For today, we'll just focus on getting the stationary distribution right, and mostly ignore runtime.

Sampling Independent Sets

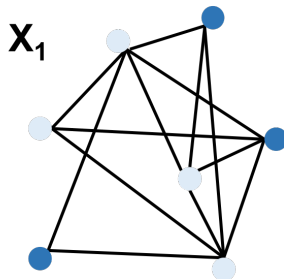
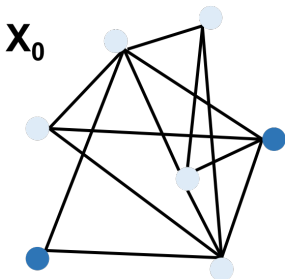
Suppose we would like to sample a uniformly random independent set from a graph G .



Very non-obvious how to sample from this distribution. Exactly counting the number of independent sets, which is closely related to sampling, is #P-hard.

Markov Chain on Independent Sets

Design a Markov chain X_0, X_1, \dots whose states are exactly the independent sets. E.g., let X_{t+1} be chosen uniformly at random from $\mathcal{N}(X_t) = \{Y : \text{independent set formed by adding/removing a node from } X_t\}$.



Unfortunately, the stationary distribution of this chain may not be uniform. It places higher probability on independent sets with lots of neighboring independent sets.

Achieving a Uniform Stationary Distribution

Define a Markov chain X_0, X_1, \dots over independent sets with transition function:

- Pick a random vertex v .
- If $v \in X_t$, set $X_{t+1} = X_t \setminus \{v\}$.
- If $v \notin X_t$ and $X_t \cup \{v\}$ is independent, set $X_{t+1} = X_t \cup \{v\}$.
- Else set $X_{t+1} = X_t$.

Is this chain irreducible and aperiodic? Yes.

For any two independent sets i, j , what is $P_{i,j}$? $P_{i,j} = P_{j,i} = 1/|V|$ if i, j differ by one vertex, $P_{i,j} = P_{j,i} = 0$ otherwise.

Thus, the Markov chain is symmetric, so by our claim from last class, the stationary distribution is uniform.

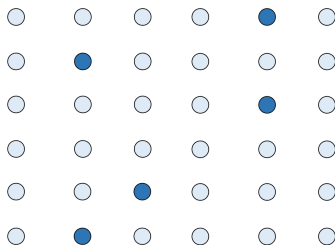
Achieving a Non-Uniform Stationary Distribution

Suppose we want to sample an independent set X from our graph with probability:

$$\pi(X) = \frac{\lambda^{|X|}}{\sum_{Y \text{ independent}} \lambda^{|Y|}},$$

for some 'fugacity' parameter $\lambda > 0$.

Known as the 'hard-core model' in statistical physics.

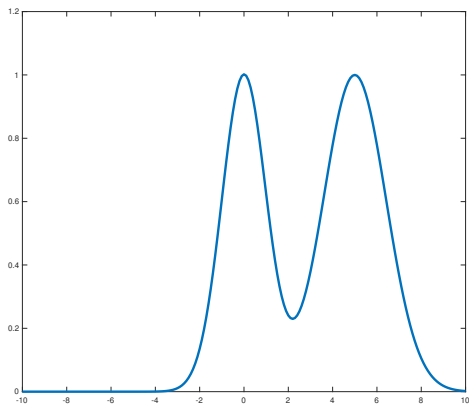


Metropolis-Hastings Algorithm

A very generic way of designing a Markov chain over state space $[m]$ with stationary distribution $\pi \in [0, 1]^m$.

- Assume the ability to efficiently compute a density $p(X) \propto \pi(X)$.
- Assume access to some **symmetric** transition function with transition probability matrix $Q \in [0, 1]^{m \times m}$.
- At step t , generate a 'candidate' state \mathbf{X}_{t+1} from \mathbf{X}_t according to Q .
- With probability $\min\left(1, \frac{p(\mathbf{X}_{t+1})}{p(\mathbf{X}_t)}\right)$, 'accept' the candidate. Else 'reject' the candidate, setting $\mathbf{X}_{t+1} = \mathbf{X}_t$.

Metropolis-Hastings Intuition



Metropolis-Hastings Analysis

Need to check that for the Metropolis-Hastings algorithm, $\pi^P = \pi$.

Suffices to show that $p^P = p$ where $p(i) \propto \pi(i)$ is our efficiently computable density.

$$\begin{aligned} [p^P](i) &= \underbrace{\sum_j p(j) \cdot Q_{j,i} \cdot \min\left(1, \frac{p(i)}{p(j)}\right)}_{\text{acceptances}} + \underbrace{p(i) \cdot \sum_j Q_{i,j} \left(1 - \min\left(1, \frac{p(j)}{p(i)}\right)\right)}_{\text{rejections}} \\ &= \sum_j Q_{i,j} \cdot \min(p(j), p(i)) + p(i) \cdot \sum_j Q_{i,j} - \sum_j Q_{i,j} \cdot \min(p(i), p(j)) \\ &= p(i) \cdot \sum_j Q_{i,j} = p(i). \end{aligned}$$

Metropolis-Hastings for the Hard-Core Model

Want to sample an independent set X with probability

$$\pi(X) = \frac{\lambda^{|X|}}{\sum_{Y \text{ independent}} \lambda^{|Y|}}.$$

- Let $p(X) = \lambda^{|X|}$ and let the transition function Q be given by:
 - Pick a random vertex v .
 - If $v \in X_t$, set $X_{t+1} = X_t \setminus \{v\}$ with probability $\min(1, 1/\lambda)$.
 - If $v \notin X_t$ and $X_t \cup \{v\}$ is independent, set $X_{t+1} = X_t \cup \{v\}$.
 - Else set $X_{t+1} = X_t$ with probability $\min(1, \lambda)$.
- Need to accept the transition with probability $\min\left(1, \frac{p(X_{t+1})}{p(X_t)}\right)$.

The key challenge then becomes to analyze the mixing time.

For the related Glauber dynamics, Luby and Vigoda showed that for graphs with maximum degree Δ , when $\lambda < \frac{2}{\Delta-2}$, the mixing time is $O(n \log n)$. But when $\lambda > \frac{c}{\Delta}$ for large enough constant c , it is NP-hard to approximately sample from the hard-core model.

Counting to Sampling Reductions

Often if one can efficiently sample from the distribution

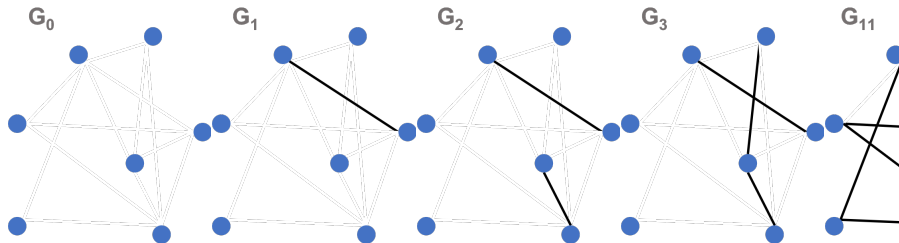
$\pi(X) = \frac{p(X)}{\sum_Y p(Y)}$, one can efficiently approximate the normalizing constant $Z = \sum_Y p(Y)$ (often called the **partition function**).

- If Z is hard to approximate, then this can give a proof that sampling is hard, and thus it is unlikely that any simple MCMC method for sampling from π mixes rapidly.
- This is e.g., how one can show that sampling from the hard-core model is hard when $\lambda = \Omega(1/\Delta)$.
- Let's consider the simple case of $\lambda = 1$. I.e., we want to sample a uniformly random independent set.
- In this case, $Z = |S(G)|$, the number of independent sets in G . It is known that approximating $|S(G)|$ even up to a $\text{poly}(n)$ factor is NP-Hard.

Counting Independent Sets

How can we count the number of independent sets $|S(G)|$ in a graph, given an oracle for sampling a uniform random independent set?

Let G_0, G_1, \dots, G_m be a sequence of graphs with $G_m = G$ and G_i obtained by removing an arbitrary edge from G_{i+1} .



We can write:

$$|S(G)| = \frac{|S(G_m)|}{|S(G_{m-1})|} \cdot \frac{|S(G_{m-1})|}{|S(G_{m-2})|} \cdot \dots \cdot \frac{|S(G_1)|}{|S(G_0)|} \cdot |S(G_0)|.$$

Counting Independent Sets

$$|S(G)| = \frac{|S(G_m)|}{|S(G_{m-1})|} \cdot \frac{|S(G_{m-1})|}{|S(G_{m-2})|} \cdots \frac{|S(G_1)|}{|S(G_0)|} \cdot |S(G_0)| 2^n = 2^n \cdot \prod_{i=1}^m r_i,$$

where $r_i = \frac{|S(G_m)|}{|S(G_{m-i})|}$. If we can estimate each r_i with \tilde{r}_i satisfying

$$\left(1 - \frac{\epsilon}{2m}\right) \cdot r_i \leq \tilde{r}_i \leq \left(1 + \frac{\epsilon}{2m}\right) \cdot r_i,$$

then:

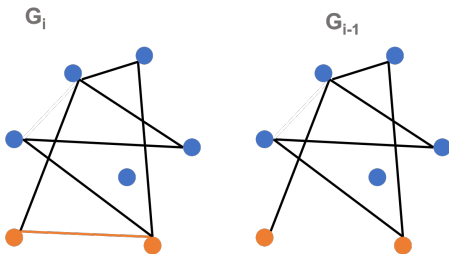
$$(1 - \epsilon) \cdot |S(G)| \leq 2^n \cdot \prod_{i=1}^m \tilde{r}_i \leq (1 + \epsilon) \cdot |S(G)|$$

since $\left(1 + \frac{\epsilon}{2m}\right)^m \leq 1 + \epsilon$ and $\left(1 - \frac{\epsilon}{2m}\right)^m \geq 1 - \epsilon$.

Independent Set Ratios

Consider the ratio $r_i = \frac{|S(G_i)|}{|S(G_{i-1})|}$. Observe that $r_i \leq 1$.

Further, $r_i \geq 1/2$. Let (u, v) be the edge removed from G_i to obtain G_{i-1} . Then each independent set in $S(G_{i-1}) \setminus S(G_i)$, must contain both u and v .



So, we can map each set in $S(G_{i-1}) \setminus S(G_i)$ to a unique set in $S(G_i)$ by simply removing v .

$$r_i = \frac{|S(G_i)|}{|S(G_{i-1})|} = \frac{|S(G_i)|}{|S(G_i)| + |S(G_{i-1}) \setminus S(G_i)|} \geq \frac{1}{2}.$$

Independent Set Ratios

So Far: We have written $|S(G)| = 2^n \cdot \prod_{i=1}^m r_i$ where $r_i = \frac{|S(G_i)|}{|S(G_{i-1})|}$.
Need to get a $1 \pm \epsilon/m$ estimate to each r_i to get a $1 \pm \epsilon$ estimate to $|S(G)|$.

Let \mathbf{X} be a random variable generated as follows: pick a random independent set from G_{i-1} and let $\mathbf{X} = 1$ if the set is also independent in G_i . Otherwise let $\mathbf{X} = 0$.

What is $\mathbb{E}[\mathbf{X}]$?

How many samples of \mathbf{X} do we need to take to obtain a $1 \pm \epsilon/m$ approximation to r_i with high probability?

Counting Independent Sets

Upshot: For a graph G with m edges, making $\tilde{O}(m^2/\epsilon^2)$ calls to a uniform random independent set sampler on G or its subgraphs suffices to approximate the number of independent sets in G up to $1 \pm \epsilon$ relative error.

- So a polynomial time algorithm for uniform random independent set sampling, would lead to a polynomial time algorithm for counting independent sets, and hence the collapse of NP to P .
- Observe that near-uniform sampling (as would be obtained e.g., with an MCMC method) would also suffice.